

# TÌM KIẾM CHUỖI

Phan Lộc Sơn 19120033

Đoàn Kim Huy

Đỗ Hoài Nam

Trần Anh Huy

Ngày 17 tháng 12 năm 2020

# TÌM KIẾM CHUỖI

## 1 Xác định vấn đề

Tìm kiếm chuỗi hiện diện rất nhiều trong cuộc sống, ví dụ như: tìm tên thầy dạy DSA trong danh sách giảng viên, tìm tên trong bảng điểm,.. hoặc trong khoa học, như là tìm liệu cấu trúc ADN của virus này có trong virus khác hay không.

Trong Tin học, các trình soạn thảo văn bản thường phải tìm kiếm (tất cả) các lần xuất hiện của một đoạn văn bản trong một văn bản dài. Thông thường văn bản được chỉnh sửa liên tục, và các phần văn bản cần tìm kiếm thì được nhập bởi người dùng. Việc phát minh ra các thuật toán tìm kiếm chuỗi hiệu quả đã giúp ích rất nhiều cho các bài toán kể trên.

Bài toán tìm kiếm chuỗi thường được mô tả theo mô hình sau:

Một chuỗi cần tìm kiếm  $S$  có dạng một chuỗi ký tự  $S[1..m]$ , cần tìm chuỗi  $S$  trong một đoạn văn bản  $T[1..n]$ , với  $m \leq n$ . Đồng thời, tất cả các ký tự trong  $T$  và  $S$  đều thuộc về một tập hữu hạn các ký tự cho trước. Chuỗi  $S$  được gọi là xuất hiện bắt đầu từ vị trí  $p + 1$  nếu thỏa điều kiện:  $T[p + j] = S[j]$ , với  $1 \leq j \leq m$ .

Có 3 thuật toán thường được dùng để tìm kiếm chuỗi:

- Thuật toán Brute-force (vét cạn, hay còn gọi là thuật trâu)
- Rabin - Karp
- Knuth - Morris - Pratt

Chi tiết của từng thuật toán sẽ được trình bày bên dưới.

## 2 Một số thuật toán tìm kiếm chuỗi

### 1. Brute-force

Giải thuật Brute-Force, hay còn gọi là vét cạn, là thuật toán đơn giản nhất trong các thuật toán tìm kiếm chuỗi con *pattern* trong chuỗi cha *text*.

Có thể giải thích đơn giản, giải thuật Brute-Force so sánh lần lượt mỗi chuỗi con *subtext* của *text* có cùng chiều dài với *pattern* với *pattern*, nếu tìm được, trả về kết quả là vị trí được tìm thấy; khi không tìm được kết quả mong muốn, trả về giá trị quy ước là không tìm thấy.

Trong ví dụ sau, ta sẽ làm rõ cách hoạt động của giải thuật này:

*text* = Let them go!

*pattern* = them

---

Let them go!

them

---

Let them go!

them

---

Let them go!  
them

---

Let them go!  
them

---

Let them go!  
them

---

Let them go!  
them

---

Let them go!  
them

---

Let them go!  
them

---

Let them go!  
them

---

Ta tìm thấy chuỗi pattern tại vị trí thứ 4!

Từ ví dụ trên, ta thiết kế mã giả cho giải thuật Brute-Force:

```
0      function: Brute-Force String Matching
1      text: chuỗi cha
2      pattern: chuỗi cần tìm
3      subtext: 1 đoạn thân của chuỗi cha có cùng chiều dài với pattern.
4      Trả về: vị trí tìm thấy đầu tiên hoặc -1 nếu không tìm thấy.
5
6      vị trí tìm thấy = -1
7      subtext = chuỗi con đầu text có độ dài bằng pattern
8      while (chưa tìm thấy hoặc chưa tới cuối text)
9          if (từng ký tự của subtext = pattern):
10             trả về vị trí
11         else:
12             dịch chuyển chuỗi con subtext trong text sang phải 1 chữ cái
13     Trả về: vị trí tìm thấy
```

Với chuỗi *pattern* có độ dài là  $M$ , chuỗi *text* có độ dài là  $N$

Phân tích độ phức tạp của thuật toán trong trường hợp xấu nhất:

- Mỗi lần so sánh với *subtext*, *pattern* phải so sánh nhiều nhất là  $M$  lần (trong trường hợp cả  $M - 1$  ký tự đầu đều đúng).
- Có tất cả  $N - M + 1$  chuỗi, vậy số chuỗi cần so sánh nhiều nhất là  $N - M + 1$  *subtext* như vậy (trong trường hợp  $N - M + 2$  chuỗi *subtext* đầu không trùng với *pattern*)  
→ Cần  $M(N - M + 1)$  lần. Vì duyệt tới cuối mảng nên đây là trường hợp tìm thấy ở cuối mảng, hoặc không tìm thấy  
→ Cận trên  $O(MN)$  (vì  $N > N - M + 1$ ).  
→ Cấp phát bộ nhớ: 0.

Phân tích độ phức tạp của thuật toán trong trường hợp tốt nhất:

- Trong trường hợp tốt nhất, có thể thấy *pattern* chính là *subtext* đầu tiên của *text*.
- Như vậy, chỉ cần tốn  $M$  lần so sánh các ký tự.  
→ Cần  $M$  lần.  
→ Cận trên  $O(M)$ .  
→ Cấp phát bộ nhớ: 0.

Độ phức tạp của thuật toán trong trường hợp trung bình:  $O(MN)$

Đánh giá:

- Dễ hiểu, thuật toán này chỉ duyệt từ đầu đến cuối, so sánh tuần tự từng chuỗi con với chuỗi cần tìm kiếm.
- Không cần bước tiền xử lý (như các thuật toán được trình bày bên dưới).
- Độ phức tạp  $O(MN)$ . Không cần xin thêm bộ nhớ.

## 2. Rabin-Karp

Thuật toán Rabin-Karp là một thuật toán được sử dụng để tìm kiếm chuỗi con *pattern* trong chuỗi cha *text* bằng cách sử dụng một hàm băm.

Hàm băm là một hàm chuyển đổi mọi chuỗi thành giá trị số, giá trị này được gọi là mã băm của nó. Ví dụ, chúng ta có thể có hàm băm  $\text{hash}(\text{"hello"})=5$ .

Giống như Thuật toán Brute-Force, thuật toán Rabin-Karp cũng dịch *pattern* qua từng phần tử trong *text* để so sánh. Nhưng sự khác biệt là thuật toán Rabin-Karp so khớp mã băm của *pattern* với mã băm của chuỗi con *subtext* của *text*, và nếu mã băm khớp thì thuật toán sẽ so sánh từng ký tự trong 2 chuỗi với nhau.

Nếu mã băm được biểu diễn bằng số nguyên không quá 64 bits, độ phức tạp thời gian (time-complexity) của việc so sánh *pattern* có độ dài  $m$  với *subtext* có cùng độ dài giảm từ  $O(m)$  xuống  $O(1)$ .

Tuy nhiên mọi thứ đều có hai mặt, vấn đề của hàm băm đó là mã băm của hai chuỗi khác nhau có thể bằng nhau. Ví dụ xét hàm băm  $\text{hash}(S)$  tính mã băm của xâu  $S$  bằng cách cộng mã ASCII của các ký tự trong  $S$ :  $\text{hash}(\text{"abcd"})=97+98+99+100=394$ , hàm băm này quá đơn giản và có khả năng gây trùng mã băm cao:  $\text{hash}(\text{"dcba"})=100+99+98+97=394$ , nhưng  $\text{"abcd"} \neq \text{"dcba"}$ .

Một hàm băm tốt thoả mãn các điều kiện sau:

- Tính toán nhanh.
- Xác suất trùng mã băm nhỏ.

Thuật toán Rabin-Karp xây dựng hàm băm với ý tưởng cơ sở: xem mọi xâu như là một chuỗi số với một cơ sở (base) nào đó. Hàm băm được tính tương tự như việc ta chuyển một số nguyên về giá trị của nó, nếu là xâu ký tự thì có thể sử dụng mã ASCII (hoặc UNICODE). Một số ví dụ:

- $\text{base}=10$ ,  $\text{hash}(\text{"425"})=4 \times 10^2 + 2 \times 10^1 + 5 \times 10^0 = 425$ .
- $\text{base}=26$ , ký tự là chữ cái từ a đến z:  $\text{hash}(\text{"abc"})=97 \times 26^2 + 98 \times 26^1 + 99 \times 26^0 = 68219$ .

Để tránh tràn số thì kết quả trên được chia lấy dư cho một số  $q$ , thường được chọn là một số nguyên tố lớn. Nếu gọi tập các ký tự được sử dụng trong chuỗi là  $\Sigma$  thì base thường được chọn sao cho  $\text{base} = |\Sigma|$  hoặc là một số nguyên tố lớn.

Độ phức tạp thời gian để tính mã băm của chuỗi độ dài  $k$  mất  $O(k)$ . Khi hiện thực thuật toán, ta sẽ "trượt" *pattern* có độ dài  $m$  trên *text* từ vị trí 0 đến  $n-m$  để so sánh mã băm. Rolling hash là hàm băm có thể tính mã băm  $h_i$  của  $\text{text}[i \dots i+m-1]$  dựa trên mã băm  $h_{i-1}$  của  $\text{text}[(i-1) \dots (i+m)]$  chỉ trong thời gian  $O(1)$  thay vì tính lại trong thời gian  $O(m)$ , từ đó tăng tính hiệu quả.

$$h_i = (\text{base} \times (h_{i-1} - \text{base}^{m-1} \times \text{text}[i-1]) + \text{text}[i+m-1]) \% q$$

Trong ví dụ sau, ta sẽ làm rõ cách hoạt động của giải thuật này:

Để đơn giản, ta chọn tập các ký tự được sử dụng trong chuỗi là  $\{a, b, c, d, e, f, g, h, i, j\}$  ứng với giá trị số lần lượt là  $\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$  và  $\text{base}=10$ ,  $q=13$ .

*text* = abccddaefg

*pattern* = cdd

Mã băm của *pattern*:  $\text{hash}(\text{"cdd"}) = (3 \times 10^2 + 4 \times 10^1 + 4 \times 10^0) \% 13 = 344 \% 13 = 6$

---

abccddaefg //  $\text{hash}(\text{"abc"}) = (1 \times 10^2 + 2 \times 10^1 + 3 \times 10^0) \% 13 = 123 \% 13 = 6$   
 cdd // Vì mã băm của *subtext* "abc" và *pattern* bằng nhau  
 // nên thuật toán so sánh từng ký tự của chúng.  
 // Do  $\text{"abc"} \neq \text{"cdd"}$  nên dịch *pattern* qua phải 1 phần tử để tiếp tục tìm kiếm.

---

abccddaefg //  $\text{hash}(\text{"bcc"}) = (10 \times (123 - 10^2 \times 1) + 3) \% 13 = 233 \% 13 = 12$   
 cdd // Vì mã băm của *subtext* "bcc" và *pattern* khác nhau  
 // nên dịch *pattern* qua phải 1 phần tử để tiếp tục tìm kiếm.

---

abccddaefg	// hash("ccd")=(10 × (233 − 10 <sup>2</sup> × 2) + 4)%13 = 334%13 = 9
cdd	// Vì mã băm của <i>subtext</i> "ccd" và <i>pattern</i> khác nhau
	// nên dịch <i>pattern</i> qua phải 1 phần tử để tiếp tục tìm kiếm.

---

abccddaefg	// hash("cdd")=(10 × (334 − 10 <sup>2</sup> × 3) + 4)%13 = 344%13 = 6
cdd	// Vì mã băm của <i>subtext</i> "ccd" và <i>pattern</i> bằng nhau
	// nên thuật toán so sánh từng ký tự của chúng.
	// Do "abc"= <i>pattern</i> nên chuỗi đã được tìm thấy.

---

Từ ví dụ trên, ta thiết kế mã giả cho giải thuật Rabin-Karp:

```

0      function: Rabin-Karp String Matching
1      text: chuỗi cha
2      pattern: chuỗi cần tìm
3      subtext: 1 đoạn thân của chuỗi cha có cùng chiều dài với pattern.
4      Trả về: vị trí tìm thấy đầu tiên hoặc -1 nếu không tìm thấy.
5
6      n=chiều dài chuỗi text
7      m=chiều dài chuỗi pattern
8      Tính mã băm của pattern và chuỗi con đầu text có độ dài bằng pattern
9      Dịch pattern từ vị trí 0, qua từng phần tử của text đến vị trí n-m
10     Kiểm tra mã băm của subtext hiện tại và pattern
11     Nếu bằng nhau, so sánh từng ký tự trong 2 chuỗi
12     Nếu hoàn toàn giống nhau, trả về vị trí
13     Tính mã băm của subtext tiếp theo, nếu âm cộng thêm một lượng q
14     Nếu trong vòng lặp không tìm thấy pattern, ta đã duyệt hết text, trả
      về -1 báo kết quả không tìm thấy pattern trong text.

```

---

Với chuỗi *pattern* có độ dài là  $M$ , chuỗi *text* có độ dài là  $N$ :

Phân tích độ phức tạp của thuật toán:

Trong quá trình tiền xử lý, để tính mã băm cho *pattern* cần duyệt qua  $M$  phần tử, tương tự với *subtext* đầu tiên *text*, vậy độ phức tạp là  $O(M + M) \in O(M)$ .

Trường hợp tốt nhất, có thể thấy *pattern* chính là *subtext* đầu tiên của *text*, khi đó chỉ cần tốn  $M$  lần so sánh các ký tự, kết hợp với quá trình tiền xử lý nên độ phức tạp là  $O(M + M) \in O(M)$ .

Trường hợp xấu nhất:

- Mỗi lần so sánh với *subtext*, *pattern* phải so sánh nhiều nhất là  $M$  lần (trong trường hợp mã băm bằng nhau và cả  $M - 1$  ký tự đầu đều đúng).
- Có tất cả  $N - M + 1$  chuỗi, vậy số chuỗi cần so sánh nhiều nhất là  $N - M + 1$  *subtext* như vậy (trong trường hợp chuỗi cần tìm nằm cuối,  $N - M + 2$  chuỗi *subtext* đầu không trùng với *pattern*).  
→ Cần nhiều nhất  $M(N - M + 1)$  lần so sánh. Vì duyệt tới cuối chuỗi nên đây là trường hợp tìm thấy ở cuối hoặc không tìm thấy và xảy ra trùng mã băm ở tất cả lần so sánh trước đó.  
→ Cận trên  $O(MN)$  (vì  $N > N - M + 1$ ).  
Trường hợp trung bình:  $O(N + M) \in O(M)$

Do trong quá trình tìm kiếm có thực hiện tính toán và lưu mã băm nên chi phí bộ nhớ là hằng số  $O(1)$ .  
Đánh giá:

- Mặc dù trong trường hợp xấu nhất độ phức tạp là  $O(MN)$  không tốt hơn giải thuật Brute-Force, nhưng giải thuật Rabin-Karp hoạt động tốt hơn nhiều trong trường hợp trung bình và thực tế.
- Có bước tiền xử lý với độ phức tạp  $O(M)$  trước khi bắt đầu tìm kiếm.
- Độ phức tạp  $O(MN)$ . Chi phí bộ nhớ  $O(1)$ .

### 3. Knuth-Morris-Pratt

Giải thuật Knuth-Morris-Pratt, về cơ bản cũng giống như thuật toán Brute-Force, tuy nhiên, chỉ khác ở chỗ, Brute-Force khi so sánh *pattern* và *subtext*, Brute-Force so sánh toàn bộ các ký tự của chúng lại từ đầu, còn Knuth-Morris-Pratt, từ lần so sánh trước, sẽ quyết định có so sánh *pattern* với *subtext* kế tiếp

không, và nếu không sẽ nhảy bao nhiêu bước đến *subtext* khác. (đã biết được chúng giống nhau), từ đó tiết kiệm được chi phí giải bài toán.

Để dễ hiểu, ta xét 1 ví dụ như sau:

```
text    = ababcababd
```

$pattern = ababd$

[illegible]

Note: Với string là chuỗi có độ dài  $n$ , và  $j$  là số nguyên dương  $0 < j < n + 1$

Chuỗi đầu: Prefix, có dạng `prefix[] = string[1, j]`, chuỗi cuối: Suffix, có dạng `Suffix[] = string[j, n]`.

Từ ví dụ trên, ta tìm cách thiết kế ý tưởng của KMP để tìm chuỗi 1 cách tổng quát:

- Ta thấy rằng, để KMP tìm cách nhảy hiệu quả, cần phải xử lý trước pattern, tìm chuỗi chung vừa là chuỗi đầu, vừa là chuỗi cuối của pattern. Bước này gọi là bước tìm Failure Function, tại mỗi vị trí  $j$  trong pattern ( $0 < j < \text{len}(\text{pattern}) + 1$ ), tìm độ dài lớn nhất của chuỗi chung và lưu vào mảng `f[]` để lưu trữ.
- Khi đang so sánh pattern và subtext, ta có thể gặp các trường hợp ( $0 < i \leq \text{len}(\text{subtext})$ ):
  1. Nếu `pattern[i] = subtext[i]`, tăng thêm  $i$  để so sánh ký tự kế tiếp.
  2. Nếu `pattern[i] ≠ subtext[i]`, và  $i = 1$ , tức là ký tự đầu của chúng khác nhau, đơn giản là chỉ cần nhảy subtext sang bên phải 1 bước.
  3. Nếu `pattern[i] ≠ subtext[i]`, và  $i > 1$ , tức là đã có 1 lượng  $i - 1$  ký tự ở trước đã được so sánh đúng, bây giờ là đến lúc cần dùng mảng `f` đã tính trước ở bên trên, ta lấy `f[i - 1]` (chính là độ dài lớn nhất của chuỗi chung của `pattern[1, i - 1]`) và ta sẽ nhảy subtext 1 lượng `min((i - 1) - f[i - 1], 1)` để chuẩn bị cho lượt so sánh kế tiếp.

Từ ý tưởng trên, ta xây dựng mã giả cho thuật toán KMP:

```
0      function: Build Failure Function
1      pattern: chuỗi cần tìm
2      Trả về: Mảng f lưu độ dài lớn nhất f[i] của chuỗi chung prefix và suffix
           của pattern[1, i].
3
4      f[] = 0 để lưu mảng trả về
5      fpre = 0;    // lưu f[i - 1] trong vòng lặp bên dưới
6      i = 2;      // Ta chỉ tính f[2] trở đi, f[1] luôn bằng 0
7      while (i <= n)
8          Nếu pattern[i] == pattern[f[i - 1] + 1]
9              fpre++
10             f[i] = fpre // tận dụng kết i - 1 kết quả so sánh của f[i - 1], ta
                chỉ cần so sánh 1 ký tự.
11             i++
12         Nếu pattern[i] != pattern[f[i - 1] + 1] và f[i - 1] = 0
13             // Tức là suffix và prefix là chuỗi có độ dài f[i - 1] + 1 có ký
                tự cuối không giống
14             // -> f[i] = 0
15             f[i] = 0
16             i++
17         Nếu pattern[i] != pattern[f[i - 1] + 1] và f[i - 1] > 0
18             // Trường hợp này, có thể f[i] = f[i - 1]
19             // Ta sẽ tận dụng pattern[1, f[i - 1]] = pattern[i - f[i - 1], i]
20             // Và f[f[j]] (độ dài này lớn nhất là j - 1)
21             // Ví dụ AAABAAAA khi xét i = 8, ta lợi dụng chuỗi [AA]ABAAAA =
                A[AA]BAAAA và [AAA]BAAAA = AAAB[AAA]A
22             // Suy ra được [AA]ABAAAA = A[AA]BAAAA = AAABA[AA]A
23             // Như thế chỉ cần xét ký tự AA[A]BAAAA và AAABAAA[A] để xác định
                f[i]
24             fpre = f[fpre]
25         Trả về: f[]
26
27      function: KMP String Matching
28      text: chuỗi cha
29      pattern: chuỗi cần tìm
30      Trả về: vị trí tìm thấy đầu tiên hoặc -1 nếu không tìm thấy.
31
32      i = 1; j = 1; (i đánh dấu vị trí đang duyệt trong text, j đánh dấu vị
                trí đang duyệt trong pattern)
33      while (i < len(text) - len(pattern))
34          Nếu pattern[j] == text[i]
35              Tiếp tục i++ và j++ và so sánh
36              Nếu j = len(pattern) trả về đã tìm thấy pattern tại vị trí i - j.
```

```

37     Nếu pattern[j] != text[i] và j == 1
38         i = i + 1 //Chuyển nhảy subtext sang phải 1 đơn vị
39     Nếu pattern[j] != text[i] và j > 1
40         j = f[i] và tiếp tục so sánh text[i] và pattern[j] // sau khi gán
           j = f[i] (giảm j) thì ta đã cố tình dịch chuyển subtext đang
           so sánh sang subtext khác!
41     Nếu trong vòng lặp không tìm thấy pattern, ta đã duyệt hết text, trả
           về -1 báo kết quả không tìm thấy pattern trong text.

```

---

Với chuỗi *pattern* có độ dài là  $M$ , chuỗi *text* có độ dài là  $N$

Phân tích độ phức tạp của thuật toán trong trường hợp xấu nhất:

- Trong quá trình tiền xử lý (tính hàm Failure Function), trong vòng while chỉ tốn nhiều nhất 1 - 2 lần để tính  $f[i]$ , vậy độ phức tạp là  $O(M)$
- Trong quá trình so sánh, mỗi lần, ta hoặc sẽ tăng  $i$  lên 1 đơn vị, hoặc sẽ giảm  $j$  ít nhất 1 đơn vị (trong trường hợp dịch  $j$ ), nghĩa là tăng  $k$  lên 1 đơn vị,  $k = i - j < N$ . Vậy tổng cộng trong vòng lặp có thể lặp  $< 2n$  lần ( $N$  lần tăng  $i$ , ít hơn  $N$  lần tăng  $k$ ), như vậy độ phức tạp là  $O(2N)$ .  
 → Độ phức tạp  $O(M + N)$ .  
 → Cấp phát bộ nhớ:  $O(M)$ .

Phân tích độ phức tạp của thuật toán trong trường hợp tốt nhất:

- Trong trường hợp tốt nhất, có thể thấy *pattern* chính là *subtext* đầu tiên của *text*.
- Như vậy, chỉ cần tốn  $M$  lần so sánh các ký tự. Ngoài ra còn tốn thêm chi phí tính Failure Function  $O(M)$ .  
 → Cận trên  $O(M)$ .  
 → Cấp phát bộ nhớ:  $O(M)$ .

Độ phức tạp của thuật toán trong trường hợp trung bình:

- Cận trên  $O(N + M)$ .
- Cấp phát bộ nhớ:  $O(M)$ .

Đánh giá:

- Nhanh, chi phí tuyến tính.
- Phải xử lý trước khi tìm kiếm thực sự.
- Độ phức tạp  $O(M + N)$ . Chi phí bộ nhớ  $O(M)$ .

# TÌM KIẾM CHUỖI TRONG TIN HỌC

Ví dụ về Crossword game

## 1 Hướng giải Crossword game

Trong phần này, nhóm chúng tôi quyết định sẽ sử dụng thuật toán KMP để làm giải thuật tìm kiếm chuỗi, bởi vì giải thuật này có vẻ nhanh hơn 2 giải thuật còn lại với thời gian  $O(M + N)$ .

Thuật toán chính của chương trình là đọc mảng chữ cái vào mảng **table** để tìm theo hàng ngang, và tạo mảng **Transpose\_table** bằng cách chuyển vị mảng **table** để tìm theo hàng dọc.

Tìm bằng thuật toán KMP, hàng ngang trước và hàng dọc sau, nếu tìm thấy, in ra vị trí tìm thấy, nếu không, in ra NF.

Độ phức tạp:

- Có  $H$  hàng, mỗi hàng cần tìm 1 chuỗi có độ dài  $M$  trong chuỗi  $W$ , vậy độ phức tạp là  $O(H * (M + W))$ .
- Có  $W$  cột, mỗi cột cần tìm 1 chuỗi có độ dài  $M$  trong chuỗi  $H$ , vậy độ phức tạp là  $O(W * (M + H))$ .  
→ Như vậy, độ phức tạp là:  $O(2WH + M(H + W))$  đối với từng từ cần tìm kiếm.

## 2 Một số mẫu kiểm tra lời giải